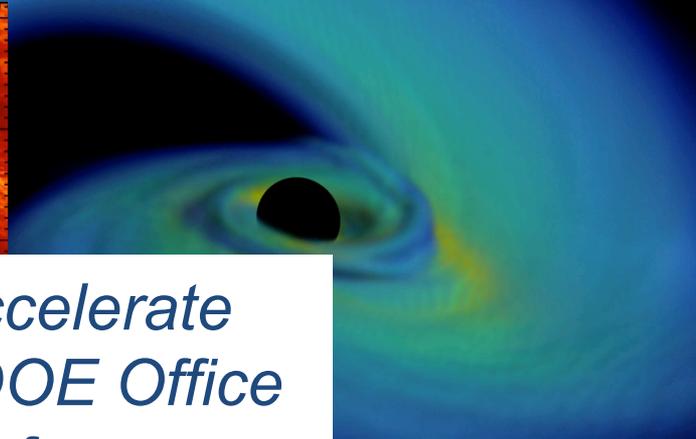
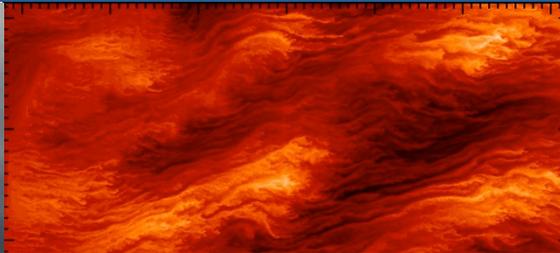


# Current and Next- Generation Supercomputing and Data Analysis at NERSC

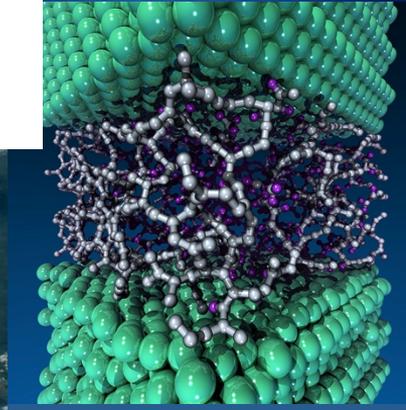
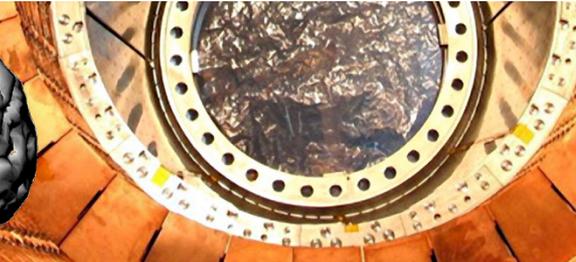
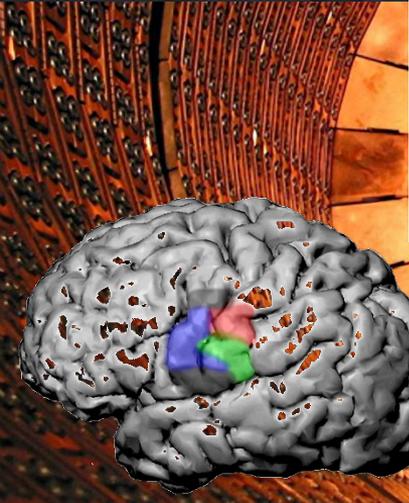


Richard Gerber

NERSC Senior Science Advisor  
High Performance Computing Department Head



*NERSC's mission is to accelerate scientific discovery at the DOE Office of Science through high performance computing and data analysis.*



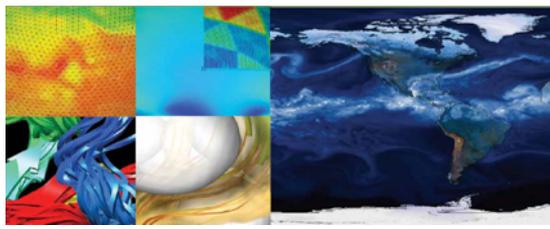
# NERSC: the Mission HPC Facility for DOE Office of Science Research



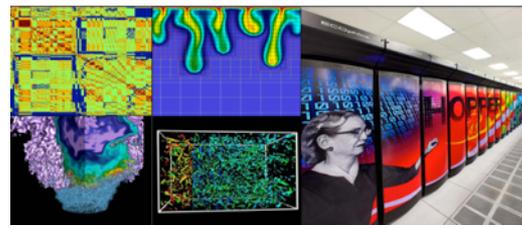
U.S. DEPARTMENT OF  
**ENERGY**

Office of  
Science

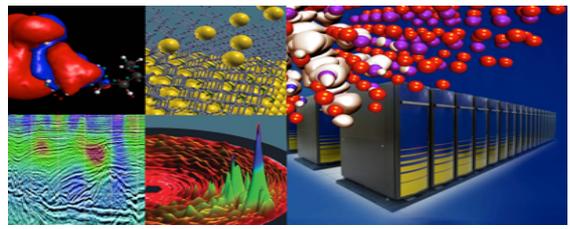
Largest funder of physical  
science research in the U.S.



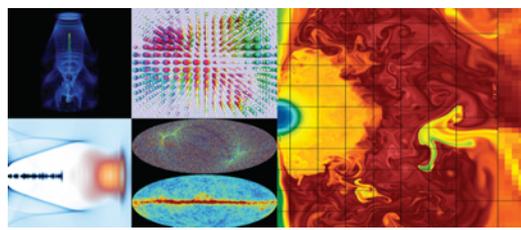
Bio Energy, Environment



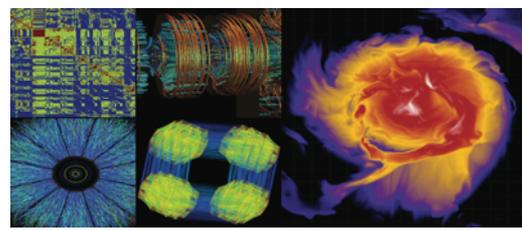
Computing



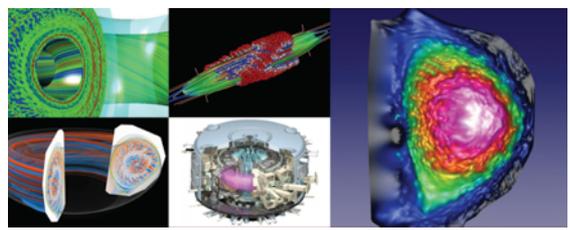
Materials, Chemistry, Geophysics



Particle Physics, Astrophysics



Nuclear Physics



Fusion Energy, Plasma Physics



U.S. DEPARTMENT OF  
**ENERGY**

Office of  
Science

- NERSC users produce more scientific publications than any other center in the world\*; ~2K/year
- 1,036 citations via Web of Science in 2017 so far (not perfect!)



5 in Nature  
30 in Nature Comm.  
70 in 12 journals

Journal	Articles
Physical Review B (condensed matter and materials)	55
Astrophysical Journal	36
Physical Review Letters	38
Journal of Chemical Physics	32
Nature Communications	30



4 in Science



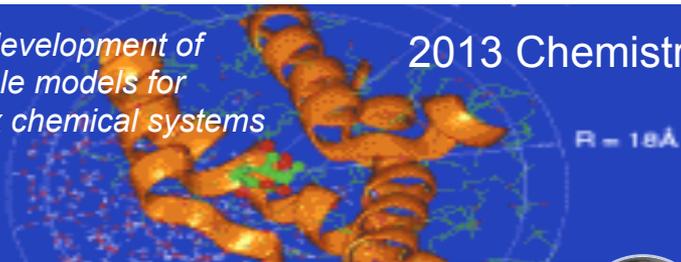
11 in PNAS

# Nobel-Prize Winning Users



*for the development of multiscale models for complex chemical systems*

2013 Chemistry

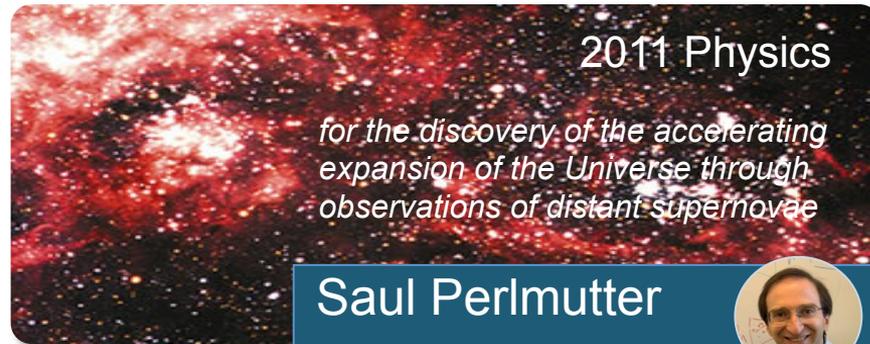


Martin Karplus



2011 Physics

*for the discovery of the accelerating expansion of the Universe through observations of distant supernovae*

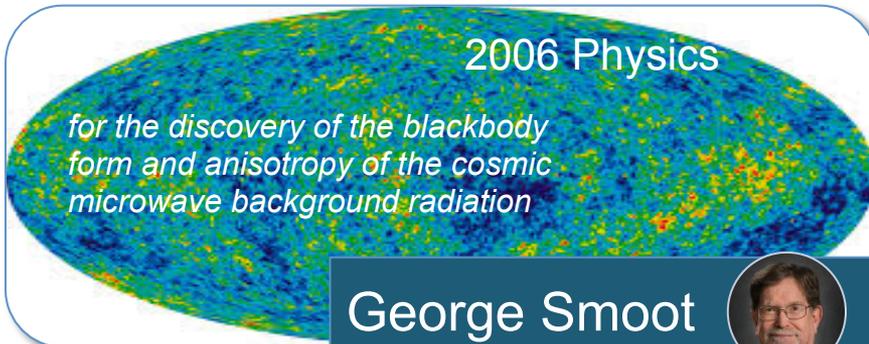


Saul Perlmutter



2006 Physics

*for the discovery of the blackbody form and anisotropy of the cosmic microwave background radiation*

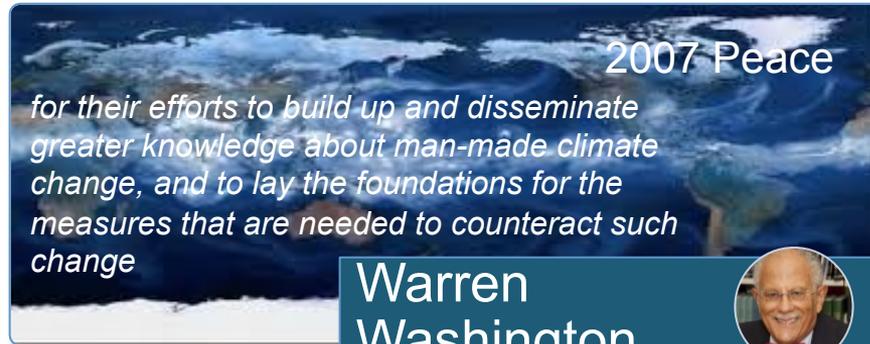


George Smoot



2007 Peace

*for their efforts to build up and disseminate greater knowledge about man-made climate change, and to lay the foundations for the measures that are needed to counteract such change*



Warren Washington



# Nobel-Prize Winning Users



*for developing cryo-electron microscopy for the high-resolution structure determination of biomolecules in solution*

2017 Chemistry

Joachim Frank

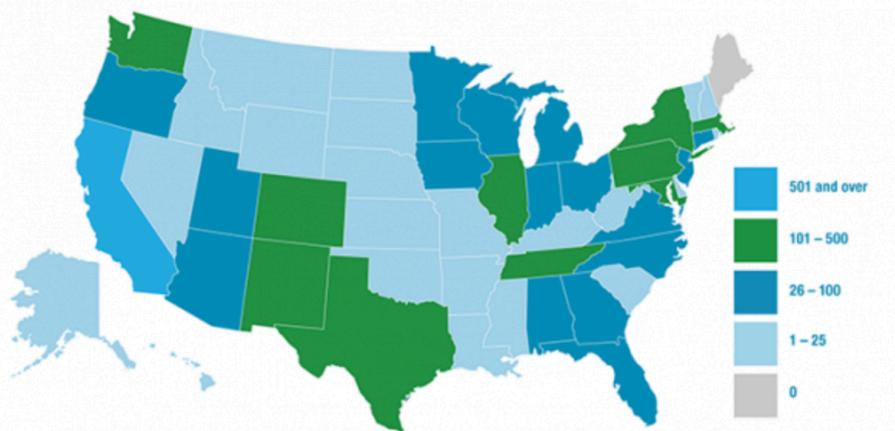
A circular portrait of Joachim Frank, a man with white hair, smiling.

*for the discovery of neutrino oscillations, which shows that neutrinos have mass*

2015 Physics

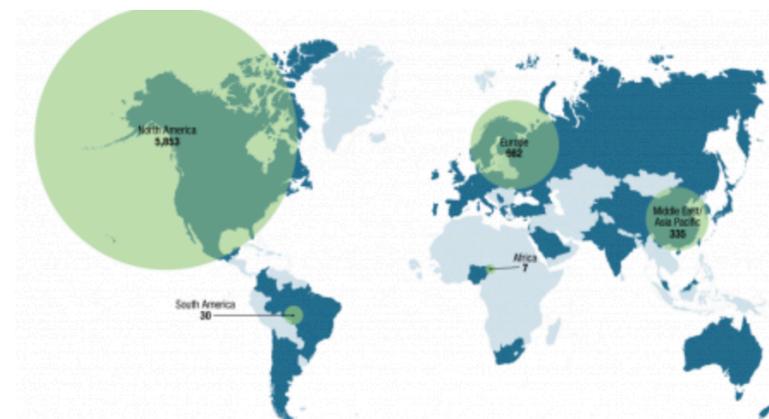
SNO Collaboration

A circular portrait of a man with glasses and a beard, likely a member of the SNO Collaboration.

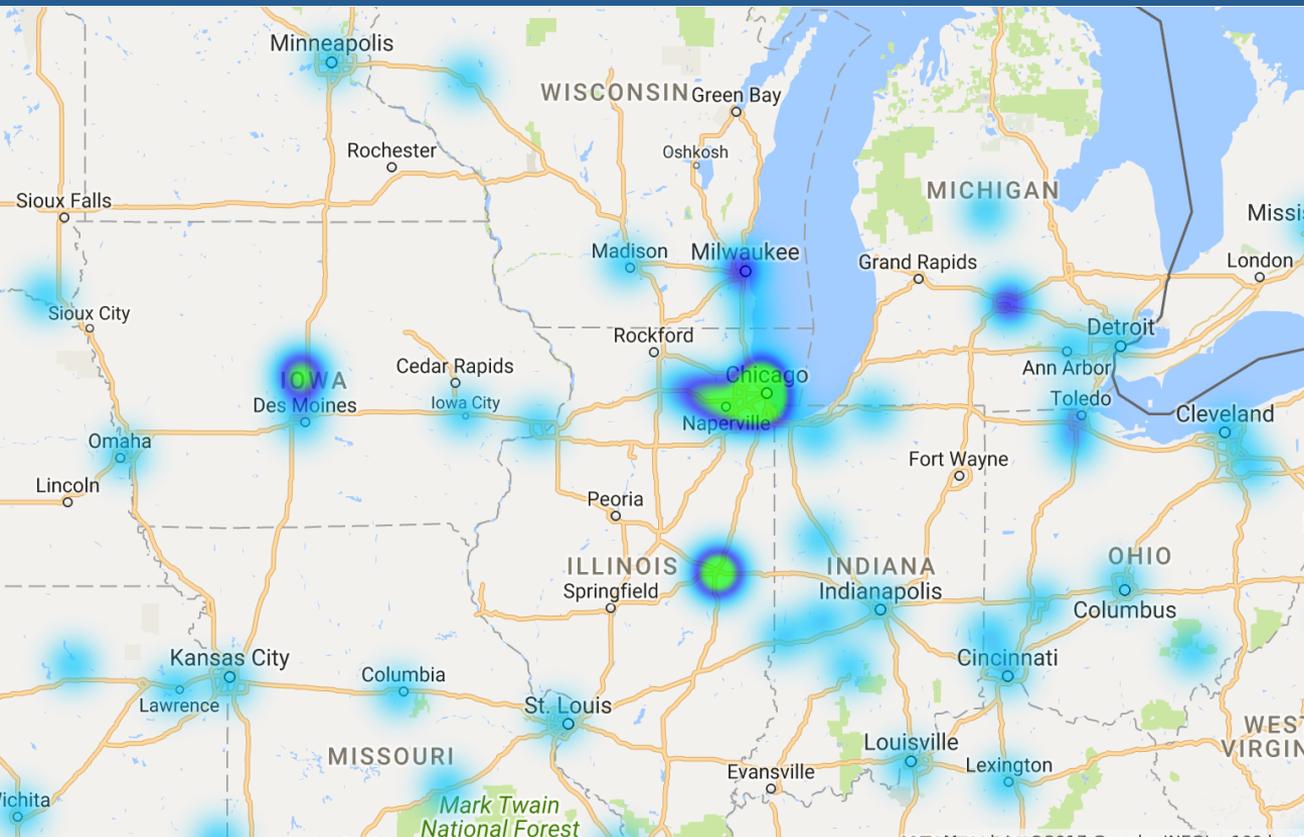


48 states  
40 countries  
Universities & national labs

7,000 users  
800 projects  
700 codes



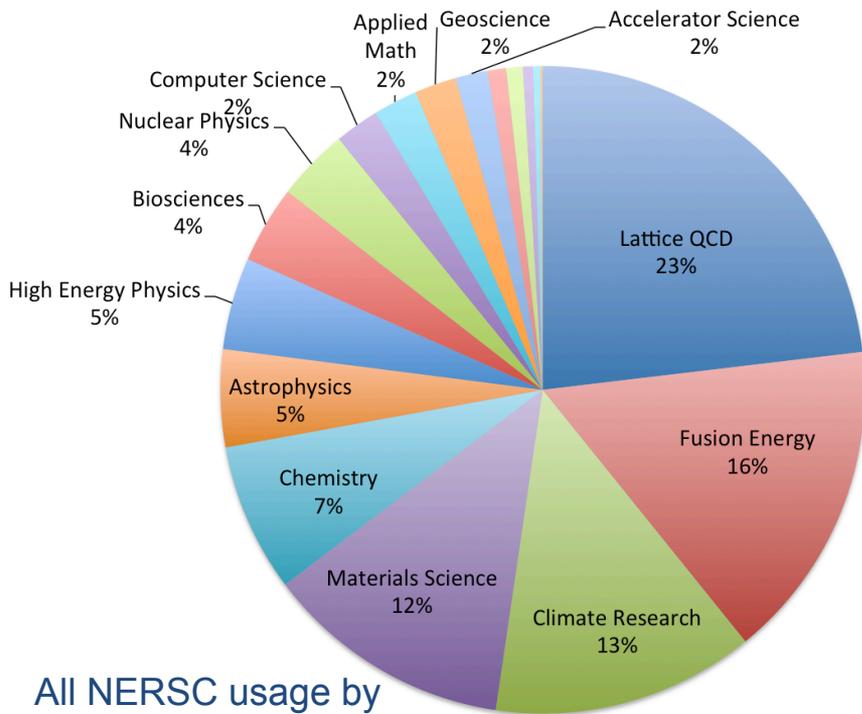
# Large contingent of active users in Iowa



63 users in Iowa

Iowa State	29
Ames Lab	26
University of Iowa	5
Drake University	1
St. Ambrose	1
Krell Institute	1

# Top Iowa Users

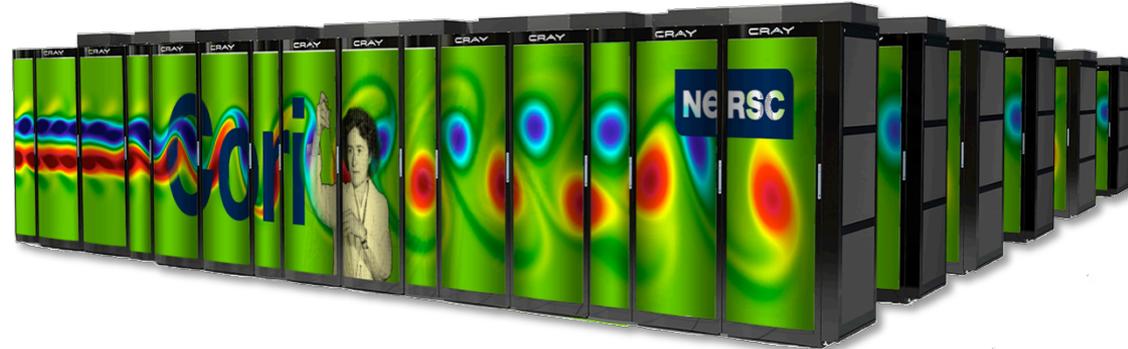


All NERSC usage by science category

Yong Han	Ames Lab	29 Mhr	Catalysis
Pieter Maris	Iowa St.	12 Mhr	Nuclear Physics
Xuewu Ou	Ames Lab	1.5 Mhr	Materials
Minsung Kim	Ames Lab	1.4 Mhr	Materials
Xin Zhao	Ames Lab	1.1 Mhr	Materials
Lin Yang	Ames Lab	0.74 Mhr	Materials
Yang Li	Iowa St.	0.53 Mhr	Nuclear Physics
F. Zahariev	Ames Lab	0.53 Mhr	Metal Extractants
Diego Floor	U. Iowa	0.50 Mhr	LQCD
M. Dick-Perez	Ames Lab	0.36 Mhr	Metal Extractants
M. C. Nguyen	Ames Lab	0.34 Mhr	Materials
Nathan Weeks	Iowa St.	0.28 Mhr	Nuclear Physics

## Cori

9,600 Intel Xeon Phi "KNL" manycore nodes  
2,000 Intel Xeon "Haswell" nodes  
700,000 processor cores, 1.2 PB memory  
Cray XC40 / Aries Dragonfly interconnect  
30 PB Lustre Cray Sonexion scratch FS  
1.5 PB Burst Buffer



#6 on list of Top 500 supercomputers in the world

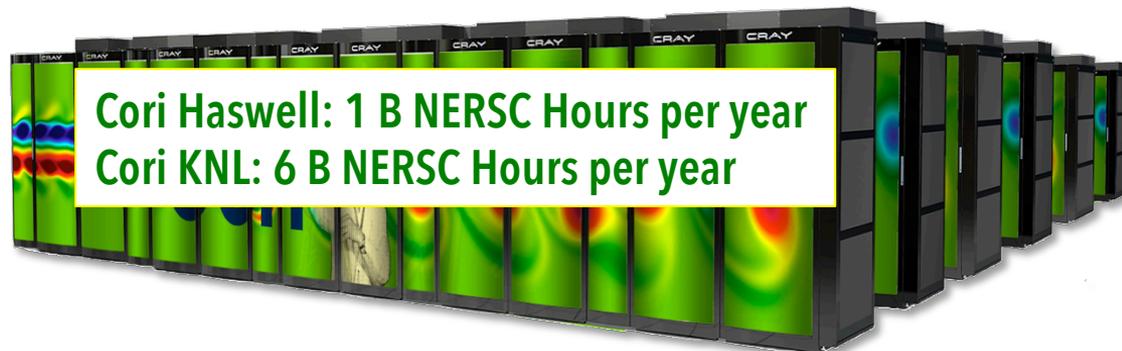


## Edison

5,560 Intel Xeon "Ivy Bridge" Nodes  
133 K cores, 357 TB memory  
Cray XC30 / Aries Dragonfly interconnect  
6 PB Lustre Cray Sonexion scratch FS

## Cori

9,600 Intel Xeon Phi "KNL" manycore nodes  
2,000 Intel Xeon "Haswell" nodes  
700,000 processor cores, 1.2 PB memory  
Cray XC40 / Aries Dragonfly interconnect  
30 PB Lustre Cray Sonexion scratch FS  
1.5 PB Burst Buffer



#6 on list of Top 500 supercomputers in the world

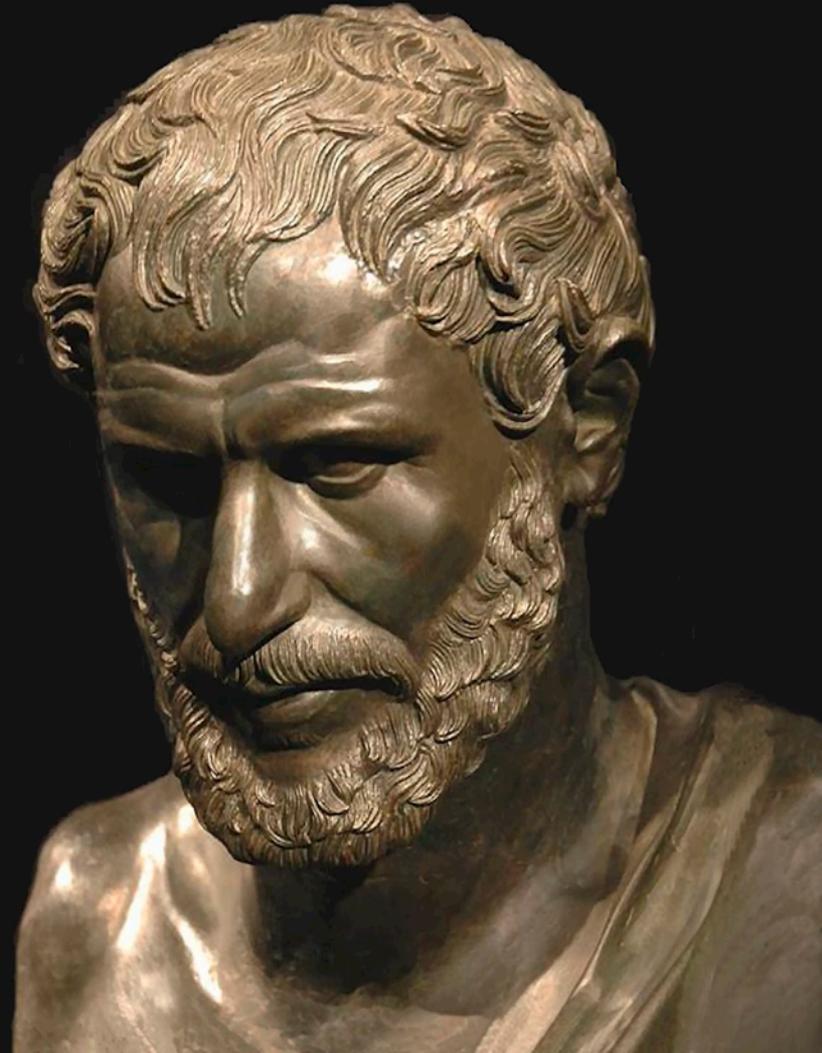


## Edison

5,560 Intel Xeon "Ivy Bridge" Nodes  
133 K cores, 357 TB memory  
Cray XC30 / Aries Dragonfly interconnect  
6 PB Lustre Cray Sonexion scratch FS

*“The only thing constant is change”*

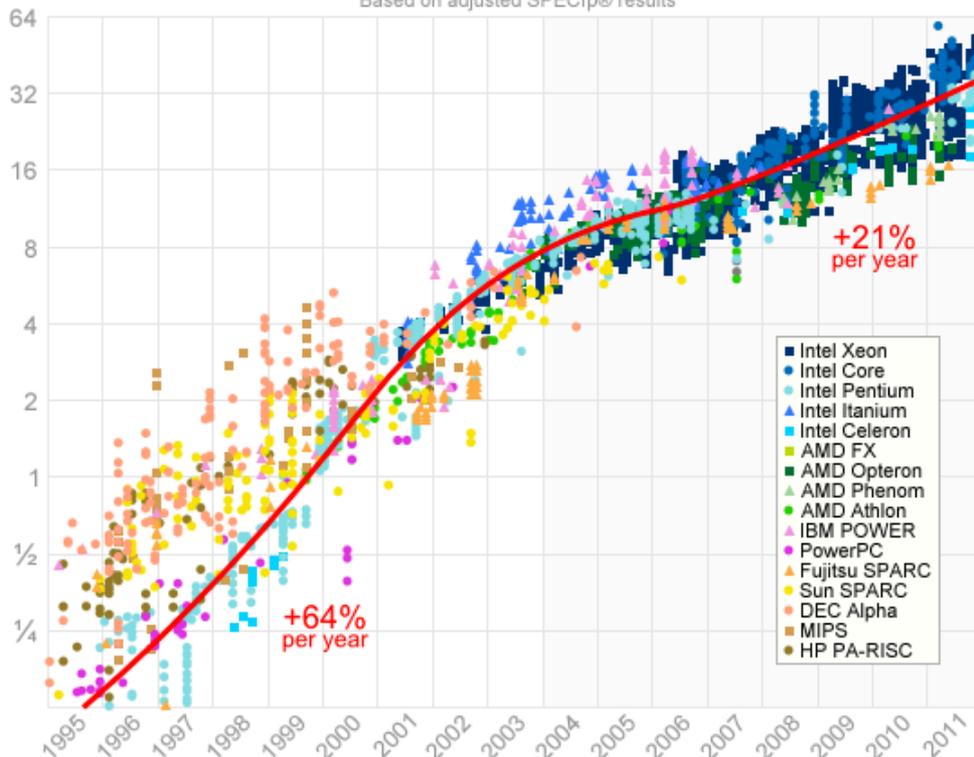
—Heraclitus of Ephesus



# Single Processor Performance

## Single-Threaded Floating-Point Performance

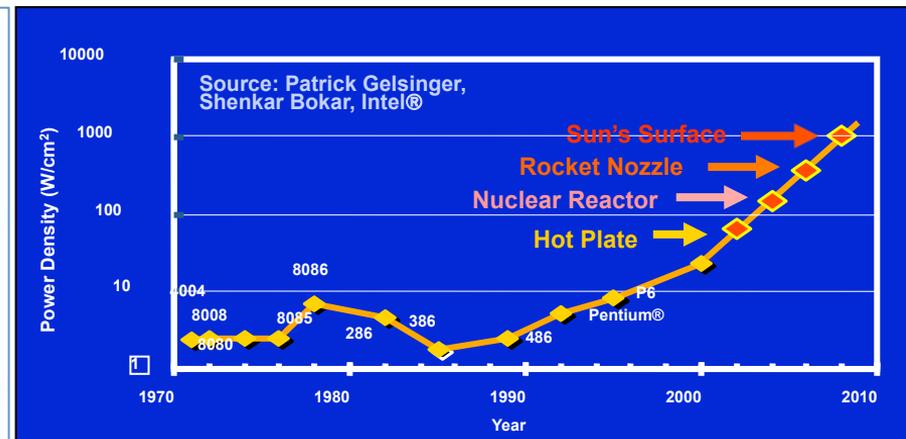
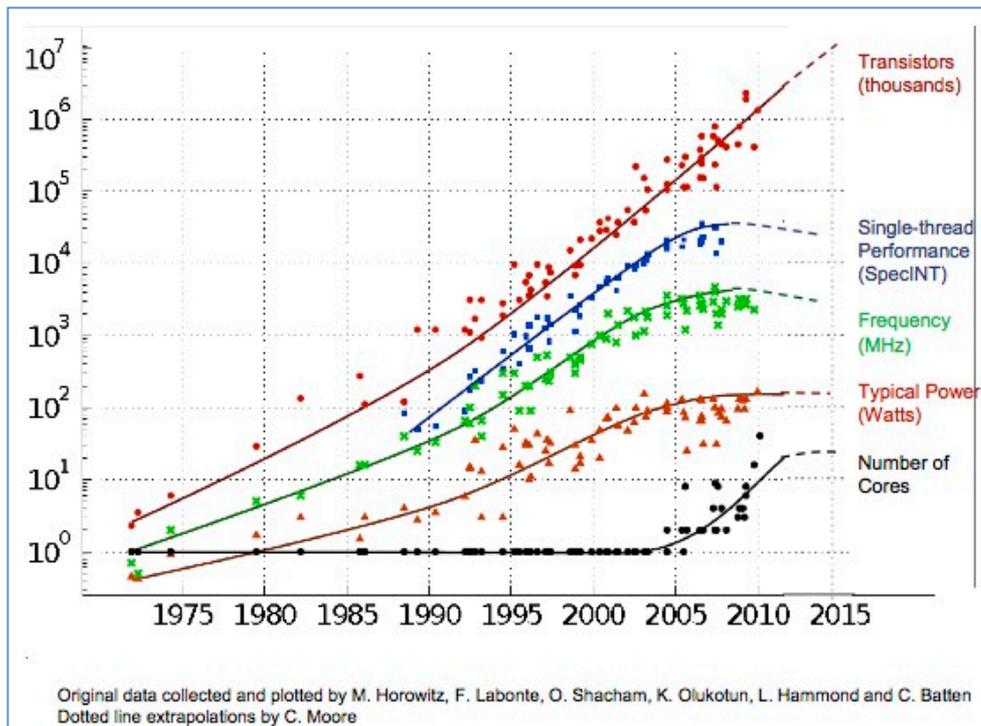
Based on adjusted SPECfp® results



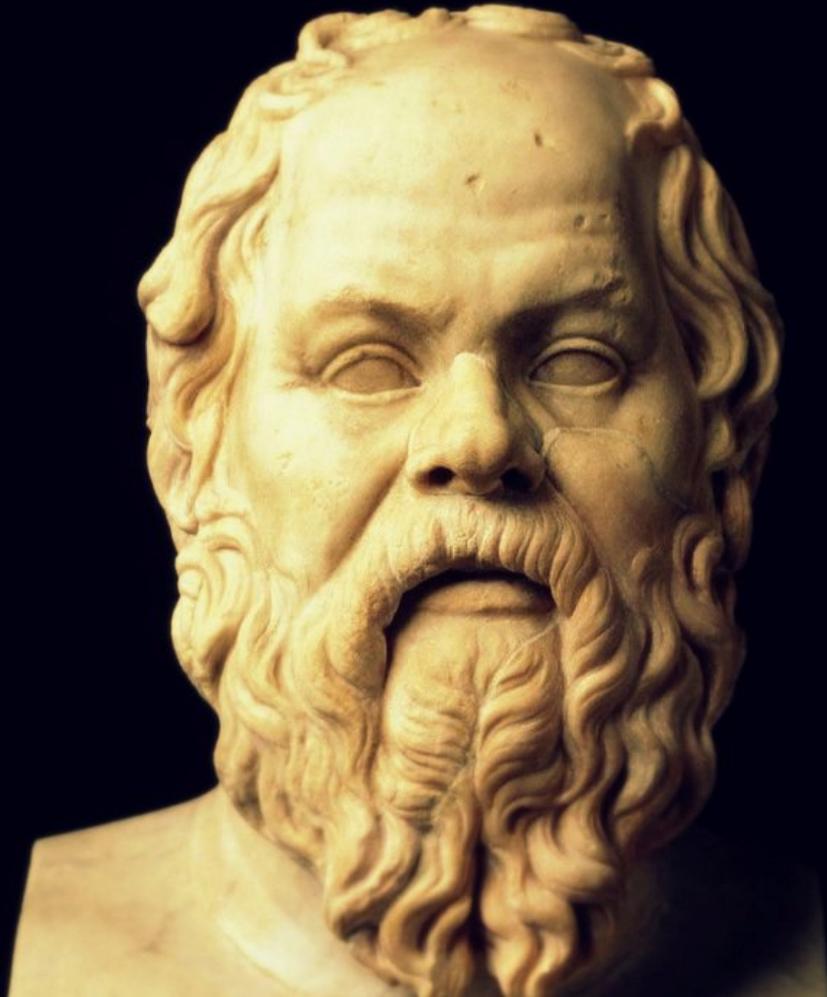
Every year there was a new CPU technology that enabled single-thread performance to increase



# Change was coming ...



Driven by power consumption and dissipation toward lightweight cores



*“The secret of change is to focus all of your energy, not on fighting the old, but on building the new.”*

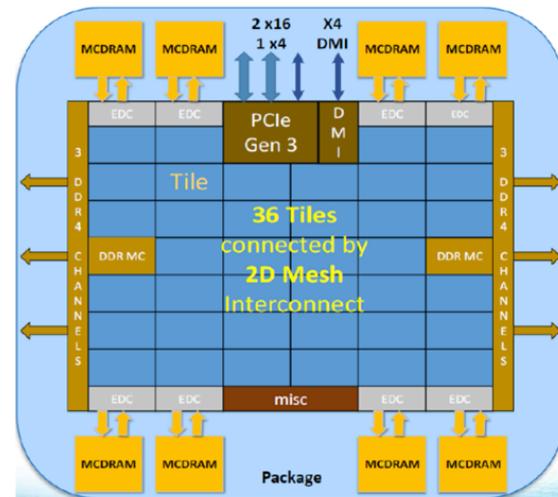
– Socrates

# NERSC to Procure "Cori" a Knights Landing Based Cray XC Supercomputer

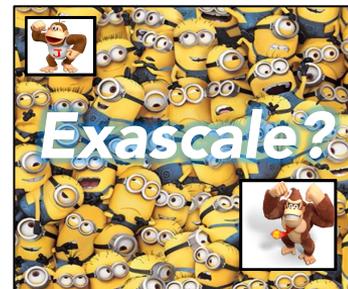
May 2, 2014 by Rob Farber — Leave a Comment

30 PFlop System will be a boon to science because of new capabilities, but the Intel Xeon Phi many-core architecture will require a code modernization effort to use efficiently.

For the first time, NERSC's users will have lower single-thread performance out of the box in their next system.



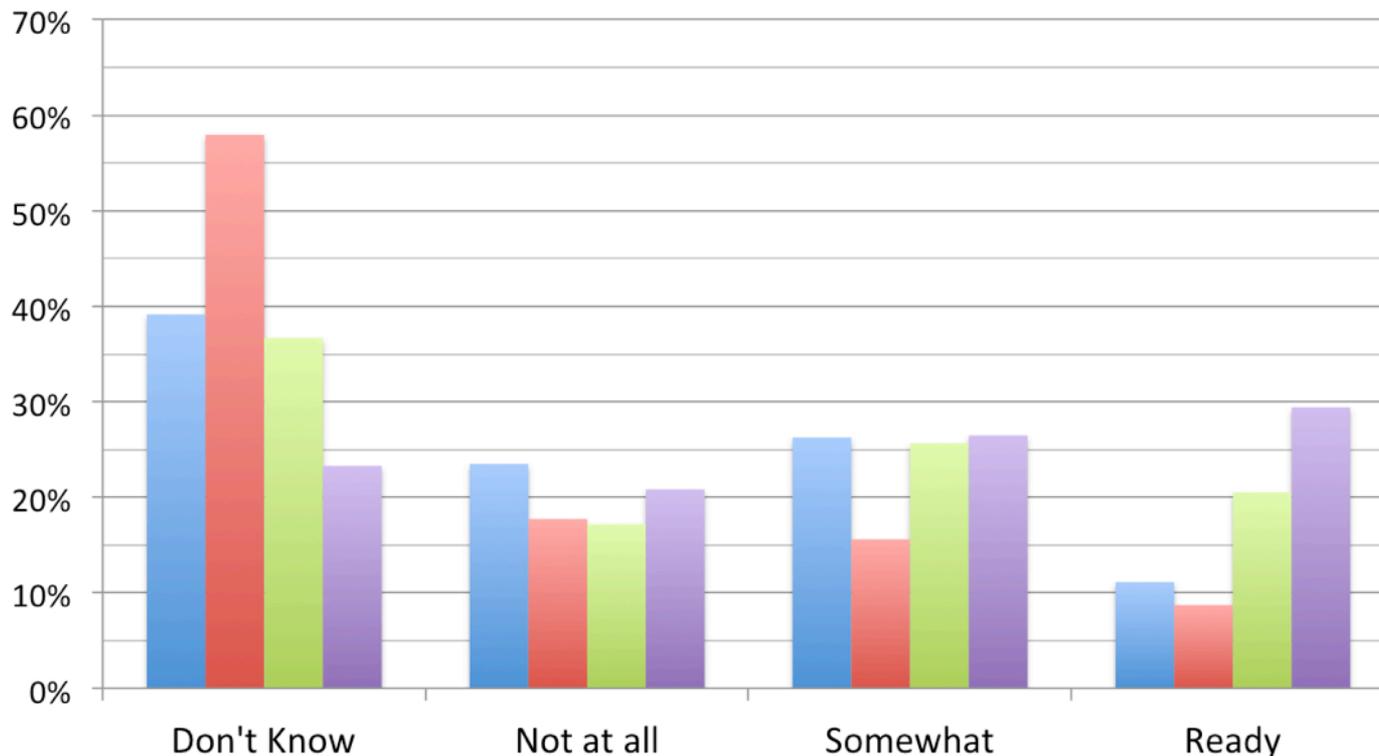
KNL: 215-230 W  
2-socket Haswell: 270 W



# 2014 User Survey: Is Your Code Ready for Manycore?



Overall    Complex memory hierarchy    Vectorization    Threading



We don't choose our users or codes. We support all DOE mission science.

Manycore is the  
future of HPC

Time to transition  
community

On the path to  
exascale

Homogeneous, x86-  
compatible CPU as a  
first step – not an  
accelerator

High bandwidth  
memory big win for  
many NERSC codes





## NERSC's Challenge

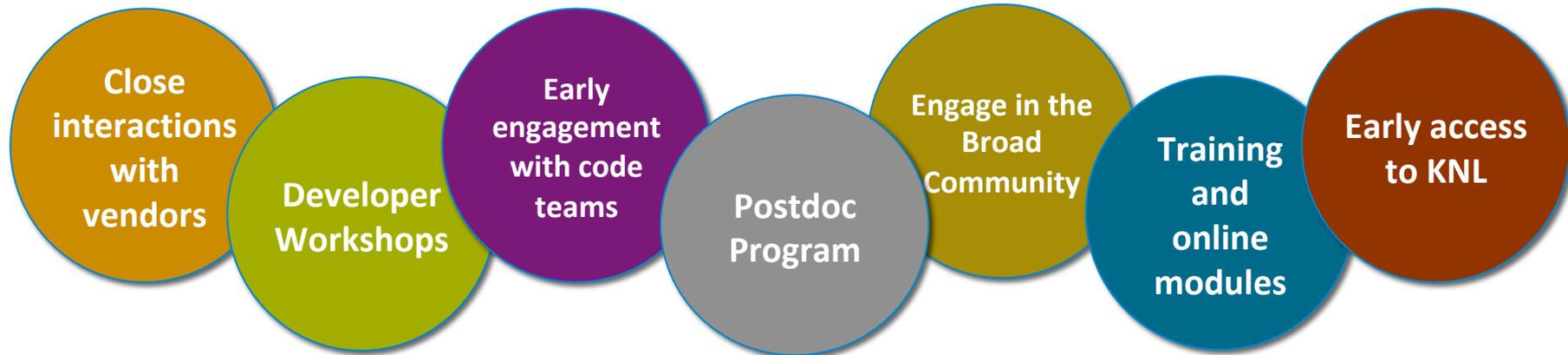
How can NERSC's diverse community of 7,000 users, 750 projects, and 700 codes use Cori's Intel Xeon Phi Knights Landing processors at high performance

Business as usual was over

# NERSC Exascale Scientific Application Program (NESAP)

---

Goal: Prepare Office of Science users for Cori's manycore CPUs  
Partner with ~20 application teams and apply lessons learned to broad user community – accounts for ~ 50% of hours used

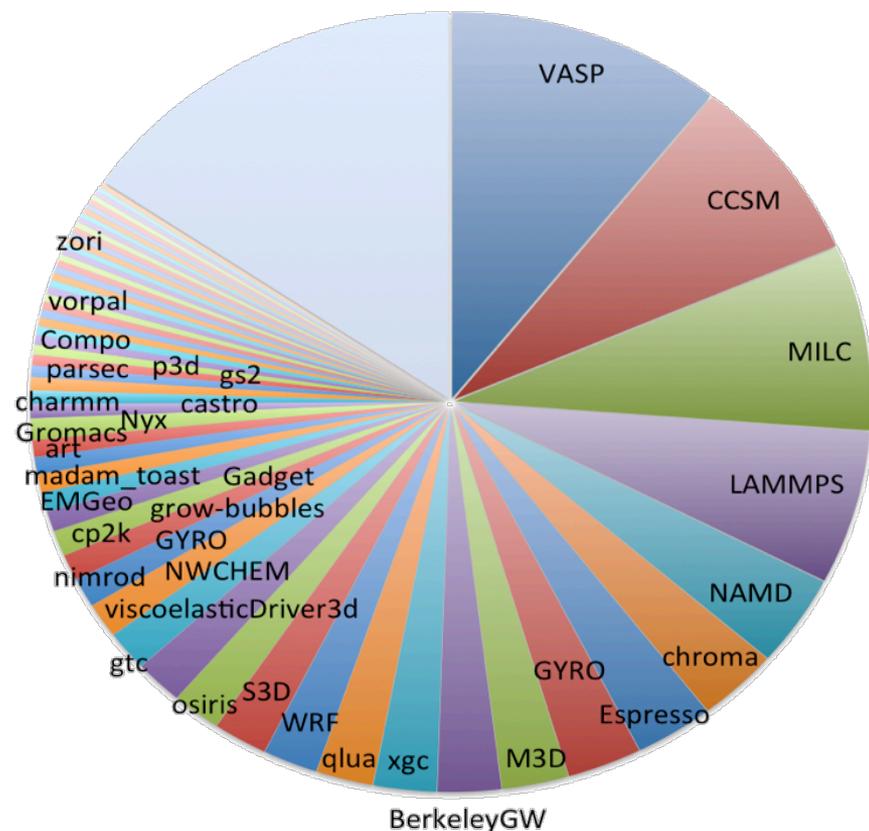


## Selected projects must

- Work with NESAP liaison to produce profiling and scaling plots and vectorization and memory BW analyses.
- **Commit 0.5-1.0 FTE to work on optimizing, refactoring, testing, and further profiling.**
- Intermediate and final reports detailing the application's science and performance improvement as a result of the collaboration.

## Evaluation criteria

- Importance to Office of Science research
- Representation all 6 OS programs
- Science potential
- Ability for code development and optimizations to be transferred to the broader community through libraries, algorithms, kernels or community codes
- Match NERSC/Vendor resources and expertise

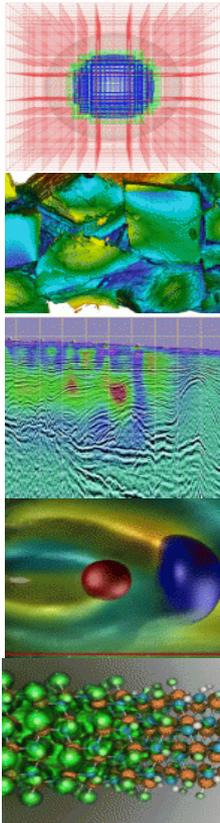


NERSC has 7,000 users, but a relatively small number of codes use a lot of hours

By working with ~20 codes, can cover ~50% of workload

A very long tail make up the last 25%

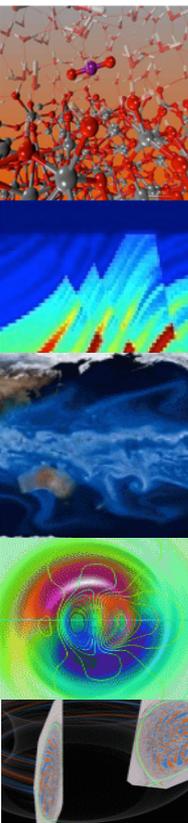
# NESAP Codes



**Advanced Scientific Computing Research**  
 Almgren (LBNL) **BoxLib** **AMR**  
 Trebotich (LBNL) **Chombo-crunch**

**High Energy Physics**  
 Vay (LBNL) **WARP & IMPACT**  
 Toussaint(Arizona) **MILC**  
 Habib (ANL) **HACC**

**Nuclear Physics**  
 Maris (Iowa St.) **MFDn**  
 Joo (JLAB) **Chroma**  
 Christ/Karsch  
 (Columbia/BNL) **DWF/HISQ**



**Basic Energy Sciences**  
 Kent (ORNL) **Quantum Espresso**  
 Deslippe (NERSC) **BerkeleyGW**  
 Chelikowsky (UT) **PARSEC**  
 Bylaska (PNNL) **NWChem**  
 Newman (LBNL) **EMGeo**

**Biological and Env Research**  
 Smith (ORNL) **Gromacs**  
 Yelick (LBNL) **Meraculous**  
 Ringler (LANL) **MPAS-O**  
 Johansen (LBNL) **ACME**  
 Dennis (NCAR) **CESM**

**Fusion Energy Sciences**  
 Jardin (PPPL) **M3D**  
 Chang (PPPL) **XGC1**

# New Postdoc Program

Open

Zahra Ronaghi  
**Tomopy**

Andrey Ovsyannikov  
**Chombo-Crunch**

Bill Arndt  
**HIPMER/  
HMMER/MPAS**

Rahul Gayatri  
**SW4**

Tuomas Koskela  
**XGC1**

Kevin Gott  
**PARSEC**

Open

**NERSC Application Performance Group formed**



Charlene Yang

# NESAP Staff Contributors



Katie Antypas



Jack Deslippe



Richard Gerber



Nick Wright



Brandon Cook



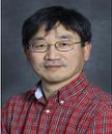
Thorsten Kurth



Helen He



Stephen Leak



Woo-Sun Yang



Rebecca Hartman-Baker



Doug Doerfler



Zhengji Zhao



Brian Austin



Rollin Thomas



Brian Friesen  
Former NESAP Postdoc

# What is different about Cori for NERSC Users?



## Edison (Cray XC w/ Intel Xeon Ivy-Bridge):

- 5000+ Nodes
- 12 Cores Per CPU
- 24 HW Threads Per CPU
- 2.4 GHz
- 8 DP Operations per Cycle
- 256b vector units
- 64 GB DDR Memory (2.6 GB/core)
- ~100 GB/s Memory BW
- 30 MB L3 cache per socket (12 cores)

## Cori (Cray XC w/ Intel Xeon Phi KNL):

- 9600+ Nodes
- 68 Physical Cores Per CPU
- 272 HW Threads Per CPU
- 1.4 GHz
- 32 DP Operations per Cycle
- 2 x 512b vector units
- 16 GB of Fast Memory (0.24 GB/core)
- 96GB of DDR Memory (1.4 GB/core)
- MCDRAM Has ~450 GB/s Memory BW

We're primarily working with existing codes to get them ready for Cori

## Goals

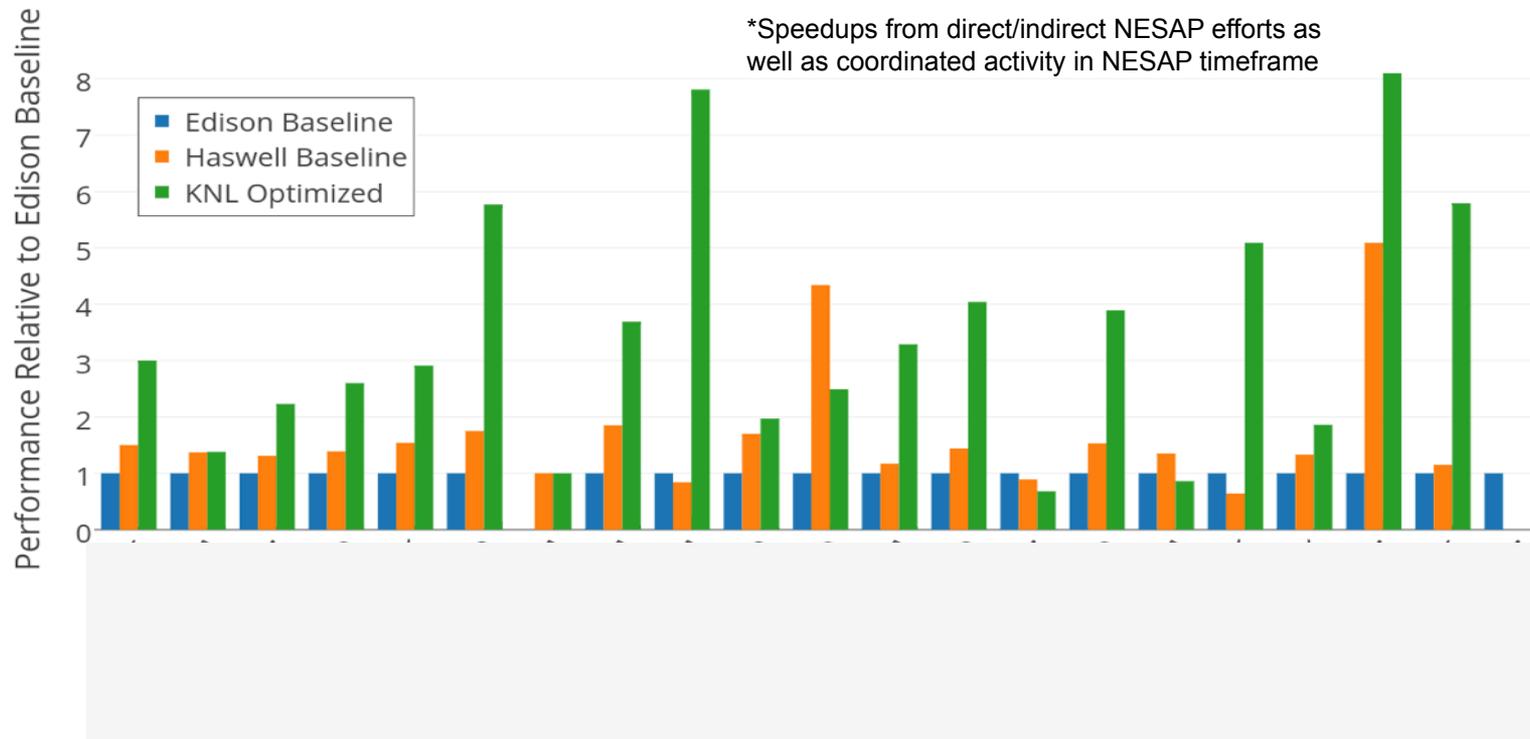
- Standard constructs for portability and maintainability
- Incorporate optimizations into code base by working directly with developers
- Collaborate closely with community to leverage expertise, communicate lessons learned, and expand NERSC influence and relevance

**Strategy:** Focus first on single-node optimization

- Enable fine-grained parallelism on light-weight cores via OpenMP
- Exploit dual 512b vector units
- Exploit 5X memory bandwidth due to MCDRAM by managing data access

Two Years Later ...

# NESAP Code Performance on KNL



# NESAP + KNL: Upgrade toward exascale



Old Code



Optimized for KNL



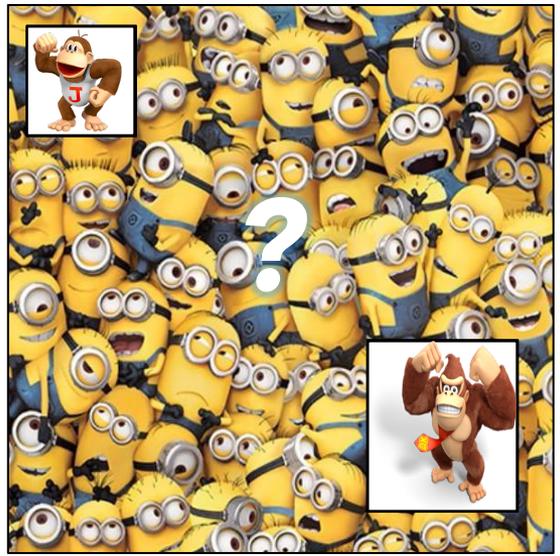
+



KNL



Exascale



Haswell



= 2.5 X performance increase per node

# Business as Usual: Dead End



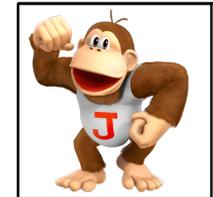
Old Code



Old Code



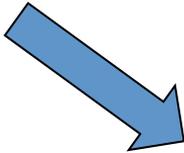
+



Edison



Haswell



1.6 X performance



Exascale

Old Code



Optimized for KNL



+



Haswell



Haswell

= 2.1 X performance

NESAP optimization efforts by themselves improved code performance by 2.1 X on x86 (Intel Xeon) processors

No motivation to optimize codes while next-gen processors gave ~60% improvements by themselves

# NESAP + KNL vs. NESAP + Haswell



Optimized for KNL



Optimized for KNL



Haswell

VS.



KNL



Exascale

1.2 X performance  
Less energy

# Remaining Challenge for the Masses



Old Code



Old Code



vs.



Haswell



KNL

= 0.7 X performance

NERSC and DOE Office of Science remaining challenge is to get broad community to run efficiently on manycore

# Summary: Good Adoption of KNL

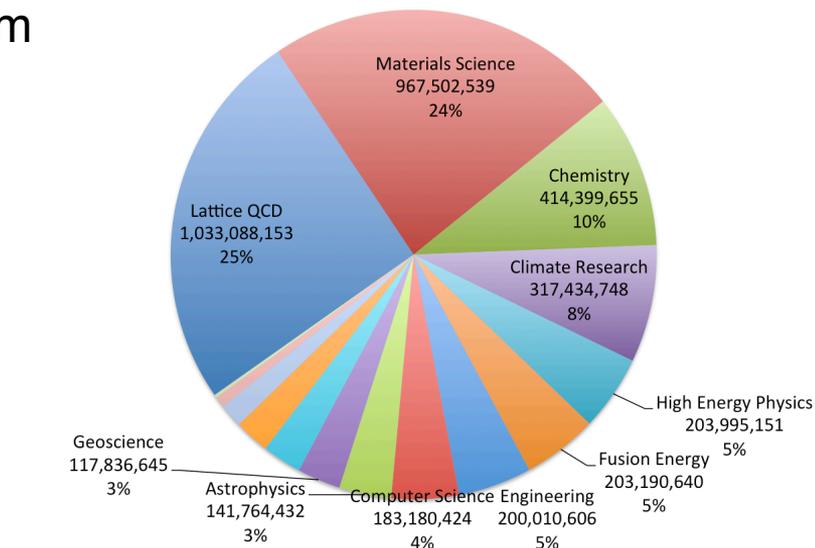


NERSC Exascale Science Application Program (NESAP) codes are running on KNL at about 3.5X-4X their pre-NESAP performance on Edison per node.

150 projects have used > 1 M NERSC Hours on KNL  
233 projects have used > 100 K NERSC Hours on KNL  
Still leaves ~500 to move over

32% of hours used by jobs using > 1,024 nodes (69K cores)

NERSC supported 6 Gordon Bell submissions using Cori KNL



Cori provides a large increase in NERSC Hours available to Office of Science researchers at NERSC (3X+ in 2018 over 2016)

# Example Science Highlights

# A Record Quantum Circuit Simulation

## Scientific Achievement

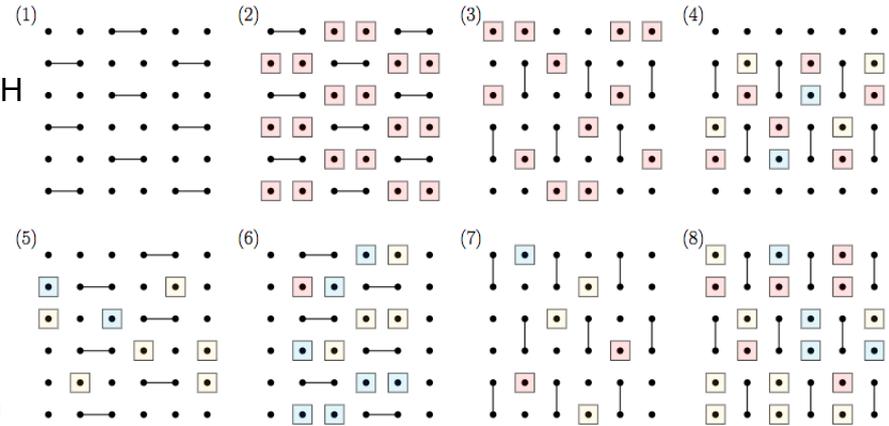
Researchers from the Swiss Federal Institute of Technology (ETH Zurich) used NERSC's 30-petaflop supercomputer, Cori, to successfully simulate a 45-qubit (quantum bit) quantum circuit, the largest simulation of a quantum computer achieved to date.

## Significance and Impact

The current consensus is that a quantum computer capable of handling 49 qubits will offer the computing power of the most powerful supercomputers in the world. This new simulation is an important step in achieving “quantum supremacy”— the point at which quantum computers finally become more powerful than ordinary computers.

## Research Details

- In addition to the 45-qubit simulation, the researchers also simulated 30-, 36- and 42-qubit quantum circuits.
- For the 45-bit simulation, they used 8,192 of 9,688 Intel Xeon Phi processors and 0.5 petabytes of memory.



Thomas Häner, Damian S. Steiger, 0.5 Petabyte Simulation of a 45-Qubit Quantum Circuit, arXiv:1704.01127 [quant-ph]

NERSC PI: T. Haner, ETH Zurich

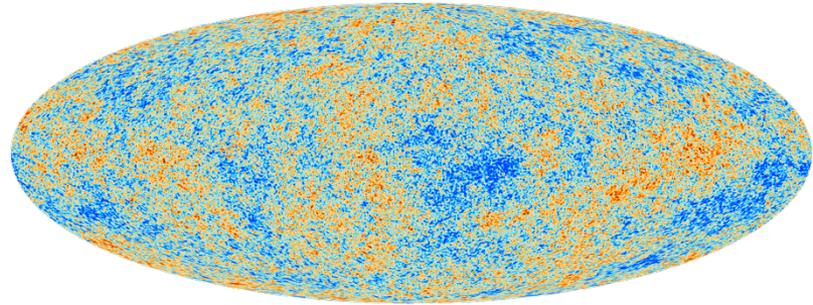
# Berkeley Lab Software Scales to 658,784 Cori Cores for Cosmic Microwave Background Analysis



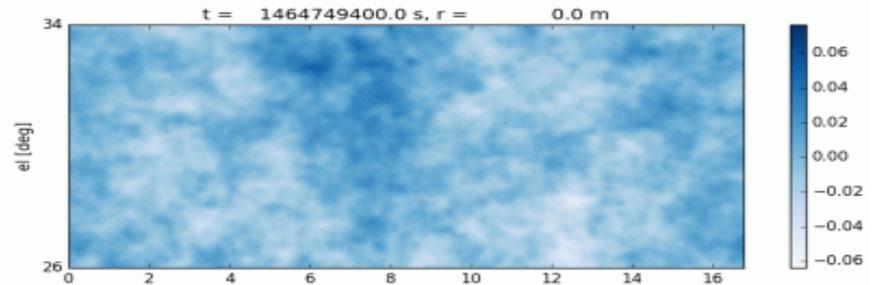
Berkeley Lab CS Computational Cosmology Center, NERSC, Intel, Cray collaborate via NESAP to scale key CMB software to full Cori system.

The TOAST (Time Ordered Astrophysics Scalable Tools) data simulation and reduction framework achieved a critical project milestone for upcoming experiments.

The ground-based CMB-S4 project will gather 35X more data than the Planck satellite did and will require TOAST's enhanced capabilities on Cori.



NERSC was used to reduce and interpret Planck results to create this map of the CMB



TOAST has incorporated modules to account for atmospheric effects



*There is no record in human  
history of a happy philosopher.*

*– H.L. Mencken*



# Data

# DOE Exascale Requirements Reviews



- Reviews with 6 Office of Science programs
- Scientists, CS researchers, facility staff
- Requirements and productivity needs for an exascale ecosystem
- Identify opportunities for collaborations among SC programs and facilities

**ASCR** ADVANCED SCIENTIFIC COMPUTING RESEARCH  
**EXASCALE REQUIREMENTS REVIEW**  
An Office of Science review sponsored by Advanced Scientific Computing Research  
SEPTEMBER 27-29, 2016  
ROCKVILLE, MARYLAND

**BER** BIOLOGICAL AND ENVIRONMENTAL RESEARCH  
**EXASCALE REQUIREMENTS REVIEW**  
An Office of Science review sponsored jointly by Advanced Scientific Computing Research and Biological and Environmental Research  
MARCH 28-31, 2016  
ROCKVILLE, MARYLAND

**BES** BASIC ENERGY SCIENCES  
**EXASCALE REQUIREMENTS REVIEW**  
An Office of Science review sponsored jointly by Advanced Scientific Computing Research and Basic Energy Sci  
NOVEMBER 3-5, 2015  
ROCKVILLE, MARYLAND

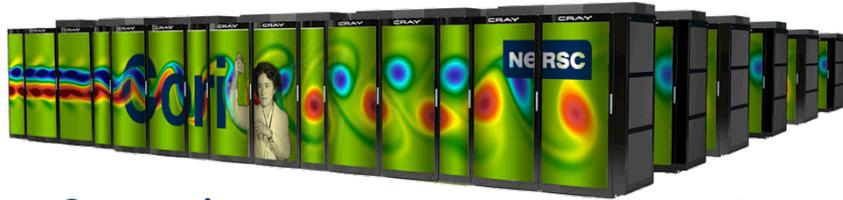
**FES** FUSION ENERGY SCIENCES  
**EXASCALE REQUIREMENTS REVIEW**  
An Office of Science review sponsored jointly by Advanced Scientific Computing Research and Fusion Energy Sciences  
JANUARY 27-29, 2016  
GAITHERSBURG, MARYLAND

**HEP** HIGH ENERGY PHYSICS  
**EXASCALE REQUIREMENTS REVIEW**  
An Office of Science review sponsored jointly by Advanced Scientific Computing Research and High Energy Physics  
JUNE 10-12, 2015  
BETHESDA, MARYLAND

**NP** NUCLEAR PHYSICS  
**EXASCALE REQUIREMENTS REVIEW**  
An Office of Science review sponsored jointly by Advanced Scientific Computing Research and Nuclear Physics  
JUNE 15-17, 2016  
GAITHERSBURG, MARYLAND



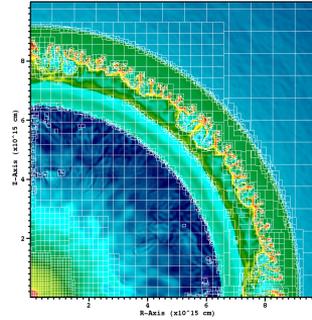
# Cross Cutting Findings



Computing



Data



Software  
and  
Applications



Workforce and Training

**CROSSCUT  
REPORT**  
**EXASCALE  
REQUIREMENTS  
REVIEWS**

March 9-10, 2017 - Tysons Corner, Virginia

An Office of Science review sponsored by:  
Advanced Scientific Computing Research  
Basic Energy Sciences  
Biological and Environmental Research  
Fusion Energy Sciences  
High Energy Physics  
Nuclear Physics

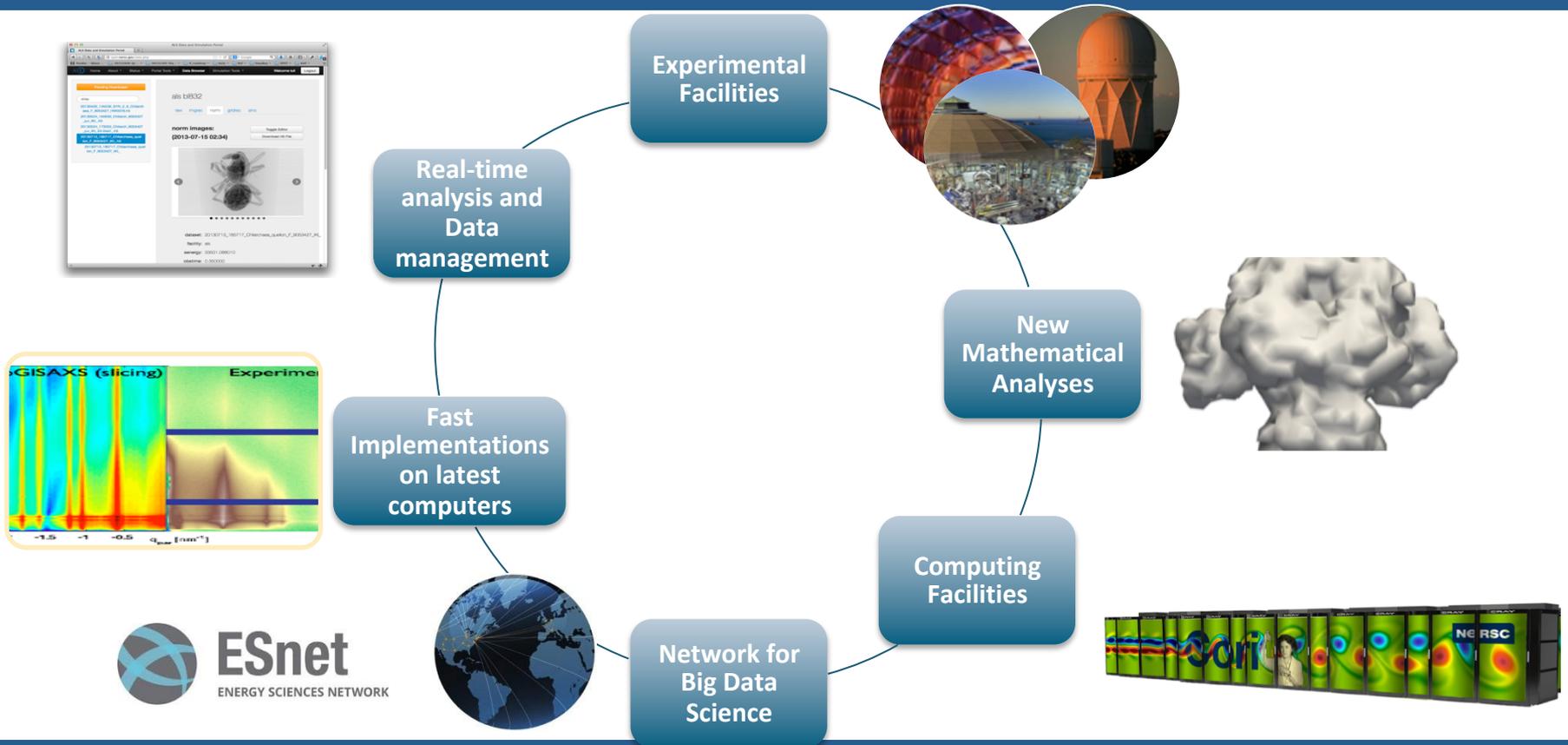
U.S. DEPARTMENT OF  
**ENERGY**

- Support for analysis **tools**, many of which **differ significantly** from traditional simulation software.
- **Complex workflows** that include data movement, co-scheduling.
- **Data management, archiving, and curation** well beyond what is in common practice today.
- Ability to **share, transfer, and access** data at remote sites with ease.
- **Input/output capability** of large HPC systems that scale with their computational capability.
- **Scheduling and allocation policies** to support workflow needs, especially data analysis requirements at experimental facilities: real-time, pseudo-real time, co-scheduling, variable job requirements, and allocations based on other resources like disk and memory.



*Transform science by enabling seamless large-scale data pipelines and analysis on leading-edge HPC systems and platforms*

# Superfacility: A network of connected facilities, software and expertise to enable new modes of discovery



# To support this vision we have created 4 initiatives, each with measurable goals



**Superfacility Initiative:** We will implement the superfacility model, collaborating with Office of Science User Facilities to enable seamless and high performing end-to-end workflows

**Systems Initiative:** We will evaluate emerging technologies and design NERSC-9 & -10 to support large-scale data analysis and simulations with tight integration to storage systems and edge services

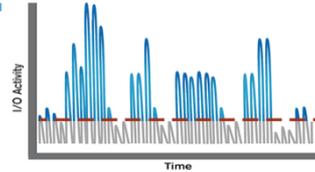
**Software Initiative:** We will partner with the community to develop and enable data analytics and management software to run at scale on HPC systems

**User Engagement Initiative:** We will engage with the user community to optimize and support data pipelines for large scale HPC systems

# NERSC has begun to address some issues raised by users



I/O is too slow



Burst Buffer more than doubles I/O bandwidth

It's difficult to bring complex software stacks to HPC systems



User defined images with Shifter



I need real-time feedback for my workflow

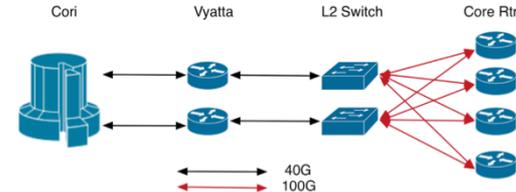


Real-time queues

Internal network limits how I can import data to supercomputer



SDN



There is limited software for analytics on HPC systems



New analytics and ML libraries

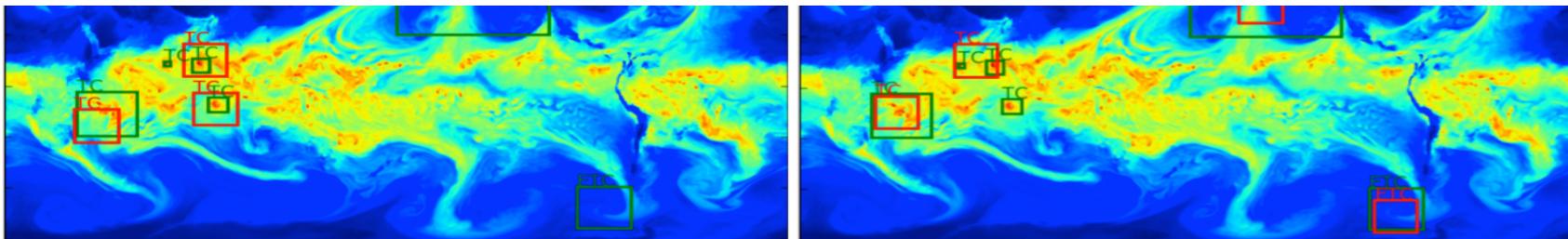


## NERSC is actively exploring Deep Learning for Science

- Collaborating with leading vendors to optimize and deploy stack
- Collaborating with leading research institutions to develop methods
- Drive real science use cases

## Deep Learning at 15 PF on NERSC Cori (Cray + Intel KNL)

- Trained in 10s of minutes on 10 terabyte datasets, millions of Images
- 9600 nodes, optimized on KNL with IntelCaffe and MKL (NERSC / Intel collaboration)
- Synch + Asynch parameter update strategy for multi-node scaling (NERSC / Stanford)



*Identified extreme climate events using supervised (left) and semisupervised (right) deep learning. Green = ground truth, Red = predictions (confidence > 0.8). [NIPS 2017]*

- Data Initiative & NESAP for Data
  - Help experimental efforts transition to KNL and towards exascale
- Continue NESAP work on Cori KNL
  - Transition broad community to manycore through training and web
  - Application portability w/ ANL, ORNL (<http://performanceportability.org>)
  - Explore 'exascale' programming models and languages
  - Influence standards committees (OpenMP, MPI)
- NESAP for NERSC 9 (2020) system when announced

# NERSC Systems Timeline



2007/2009	NERSC-5	Franklin	Cray XT4	102/352 TF
2010	NERSC-6	Hopper	Cray XE6	1.28 PF
2014	NERSC-7	Edison	Cray XC30	2.57 PF
2016	NERSC-8	Cori	Cray XC	30 PF
2020	NERSC-9			100PF-300PF
2024	NERSC-10			1EF

